

# Sniffer: Origin Tracing and Detecting of LLMs

Anonymous ARR submission

## Abstract

The extraordinary performance of large language models (LLMs) heightens the importance of detecting whether the context is generated by an AI system. More importantly, while more and more companies and institutions release their LLMs, the origin can be hard to trace. Since LLMs are heading towards the time of AGI, similar to the origin tracing in anthropology, it is of great importance to trace the origin of LLMs. In this paper, we first raise the concern of the origin tracing of LLMs and propose an effective method to trace and detect AI-generated contexts. We introduce a novel algorithm that leverages the contrastive features between LLMs and extracts model-wise features to trace the text origins. Our proposed method works under both white-box and black-box settings therefore can be widely generalized to detect various LLMs.(e.g. can be generalized to detect GPT-3 models without the GPT-3 models). We construct extensive experiments to examine whether we can trace the origins of given texts. We provide valuable observations based on the experimental results, such as the difficulty level of AI origin tracing, and the AI origin similarities, and call for ethical concerns of LLM providers. We are releasing all codes and data as a toolkit and benchmark for future AI origin tracing and detecting studies. <sup>1</sup>

## 1 Introduction

Using LLMs such as ChatGPT and GPT4 (OpenAI, 2023) for various daily routines and copilot in work is becoming a new trend that draws worldwide attention not only in the machine learning community. Starting from GPT (Radford et al., 2018) and BERT (Devlin et al., 2018), pre-trained models have developed for several years, and trustworthy and security concerns have been constantly discussed (Bai et al., 2022). The

<sup>1</sup>We are releasing all available resource at <https://github.com//>.

performances of LLMs are sensational, therefore, the usage of LLMs should be strictly supervised by users as well as service providers. One trend in ensuring the safety of LLMs is to build detection tools that can discriminate whether an AI system generates a certain text (Jawahar et al., 2020). AI-generated context detection is useful in releasing texts that require strict censoring or originality such as official documents, consultation, and student submissions (Mitchell et al., 2023) to avoid abuse of AI systems. Further, a more critical and applicable field is to trace the origin of LLMs. While more and more companies and institutions are releasing their original LLMs, it is of great importance to trace whether an LLM is trained from a previous model, or is copied or distilled from another LLM. Since LLMs can produce massive generated data, future LLMs might be trained from these generated data, the human-written texts might be contaminated with different LLMs. Therefore, tracing the origin of the text is a major challenge in future LLM industries.

In this work, we first introduce the concept of origin tracing. Then we provide an effective tool named **Sniffer** and its evaluation benchmark to study the origin tracing problem. Origin tracing is to further categorize the origin of a given text. Specifically, we categorize the origins by the LLM abilities and by LLM service providers. We first trace the origin of GPT2 level LLMs such as GPT2 (Radford et al., 2019) from OpenAI, GPT-Neo/J (Black et al., 2021; Wang, 2021) from EleutherAI, we can also trace the GPT3 level LLMs such as text-davinci-003, ChatGPT (GPT3.5 turbo) (Brown et al., 2020), LLaMA (Touvron et al., 2023). With origin tracing, we can avoid AI abuse or potential model theft since the texts can be traced to a specific service

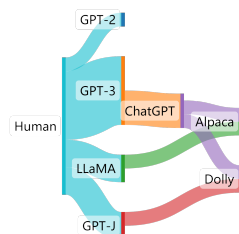


Figure 1: A knowledge flow of LLMs: with

082 provider. While more  
083 and more companies  
084 and institutions are  
085 releasing their LLMs,  
086 origin tracing is the  
087 keystone of *anthropology*  
088 of LLMs. Then we  
089 introduce **Sniffer**, the  
090 first origin tracing tool.  
091 In Sniffer, we use contrastive features across  
092 open-source LLMs such as GPT2, GPT-Neo/J, and  
093 LLaMA. Specifically, we design heuristic features  
094 that capture the model-wise discrepancies which  
095 can help trace the origin of given texts. Then  
096 we utilize a simple linear classifier to project the  
097 extracted features to specific origins including  
098 known origins that have open-source models and  
099 unknown origins that are black boxes to users.  
100 The core motivation of the origin tracing tool  
101 Sniffer utilizes the discrepancies between LLMs  
102 as features to help trace the origins. Compared  
103 with previous methods, Sniffer utilizes model-wise  
104 features, which is different from supervised  
105 learning methods such as fine-tuning a RoBERTa  
106 (Liu et al., 2019) model; Furthermore, Sniffer  
107 is able to trace text origins and generalize to  
108 unknown origins, while previous methods that use  
109 model-wise features cannot.

110 With Sniffer, we further introduce a test dataset  
111 that contains collected texts from different origins  
112 as a benchmark to study the origin tracing problem.  
113 We provide plenty of experiments and through the  
114 experimental results, we have several non-trivial  
115 observations that help future LLM studies. In  
116 general, we find that: (1) we are able to trace the  
117 origins of generated texts when we can possess the  
118 models; (2) it grows harder to detect and trace  
119 origins when the LLMs are stronger; (3) LLM  
120 providers need to be more cautious as the origin  
121 of generated texts might be known only by the  
122 providers. (4) We can trace the origins of distilled  
123 LLMs such as Alpaca and Dolly.

124 To summarize, in this paper, we: (1) raise the  
125 concern of origin tracing of AI-generated contexts;  
126 (2) build a tool to trace the origins of various open-  
127 source models by releasing a diversified benchmark  
128 for AI-generated contexts detection and origin  
129 tracing; (3) conduct various experiments to analyze  
130 the ability of LLMs when they are being traced and  
131 hope that future works can pay more attentions to  
132 the origin tracing of LLMs.

## 2 Related Work 133

### Increasing Concern of LLM Security 134

135 Pre-trained models (Devlin et al., 2018; Radford  
136 et al., 2018, 2019; Raffel et al., 2020; OpenAI,  
137 2023) are supposed to be harmless, harmful, and  
138 honest to users (Bai et al., 2022), however, there  
139 are various aspects that challenge LLM securities  
140 such as social bias, stereotypes, privacy leak or  
141 adversarial examples (Zhao et al., 2017; Carlini  
142 et al., 2020; Li et al., 2020).

143 Detection of AI-generated contexts is also a  
144 rapidly growing field that requires attention since  
145 the abuse of AI might be a major challenge in  
146 LLM applications (Bai et al., 2022). Still, current  
147 methods only consider detection as a binary task,  
148 that is, whether a given text is generated by an  
149 AI without considering tracing the origins of texts.  
150 Detection methods can be categorized into two  
151 lines:

#### Semantic-wise Detection 152

153 The most straightforward method to detect  
154 AI-generated contexts is to construct a text  
155 classification task (Zellers et al., 2019). Therefore,  
156 DetectGPT (Mitchell et al., 2023) introduces a  
157 strong RoBERTa-trained baseline that uses the  
158 dataset released by OpenAI<sup>2</sup> to train a classifier.

159 **Model-wise Detection** Unlike semantic-wise  
160 detection which discriminates the semantic  
161 difference between human-written and AI-  
162 generated texts, a more direct way is to explore  
163 the model-wise features. That is, AI-generated  
164 texts show discrepancies compared with humans  
165 when fed into AI models. These discrepancies are  
166 not easily noticed by humans as they might be  
167 a subtle difference that is possessed by a certain  
168 origin, therefore, many works focus on utilizing  
169 model-wise features such as log-likelihood of  
170 model outputs, neural features, or bag-of-word  
171 features (Bakhtin et al., 2019; Solaiman et al.,  
172 2019; Gehrmann et al., 2019; Jawahar et al.,  
173 2020; Mitchell et al., 2023) to detect AI-generated  
174 contexts.

175 In the era of LLMs, there are similar works  
176 that work on using model-wise features of LLMs  
177 such as watermarking (Kirchenbauer et al., 2023),  
178 and backdoor plantings (Kurita et al., 2020; Li  
179 et al., 2021). In the computer vision field, model-  
180 wise features are also used in fake image detection  
181 (Dolhansky et al., 2020) but are less related to

<sup>2</sup><https://github.com/openai/gpt-2-output-dataset>

Table 1: Different Detection Tools comparisons. Generalization ability represents how well white-box methods can detect texts generated from unknown models, therefore the supervised learning methods are not applicable in the generalization ability comparison.

Method	Model Access	Data Request	Generalization	Origin Tracing
$\log p(x)$	White-Box	Zero-Shot	Low	✗
DetectGPT	White-Box	Zero-Shot	Low	✗
Supervised-Learning	Black-Box	Full Data	N/A	✗
Sniffer	Black & White Box	Low-Resource	High	✓

LLM origin tracing due to the continuous nature of images.

### 3 Methods

We aim to detect whether a context is generated by an LLM system and trace the origin of the texts. Therefore, we design a simple method **Sniffer**<sup>3</sup> that is applicable in both white-box and black-box settings and only requires limited supervised data.

The core idea is to utilize the contrastive features between different accessible language models such as GPT-2, GPT-Neo, GPT-J, and LLaMA. We first obtain the perplexity of a target sample  $S$  based on different models denoted as  $\theta_0, \theta_1, \dots, \theta_N$ , then we craft several heuristic features and construct a simple linear classifier to classify the origin of the given sample. Through such a feature engineering process, we can trace the origin of the target sample down to a known model  $\theta_n$ , we can also generalize the contrastive features to show differences between unknown source models or human-written texts.

Compared with previous detection methods, as seen in Table 1, our proposed method is the first to allow origin tracing. Model-wise detection methods such as DetectGPT are designed to detect whether a text is generated by a certain model, which can not be well generalized to origin tracing.

#### 3.1 Steps of Sniffer

The process of Sniffer includes: (1) Obtain and align token-level perplexity between different models; (2) Extract contrastive features; (3) Train features for origin tracing.

##### Obtain and Align Single Model Perplexity

Given target text  $S$  and a known model  $\theta_n$ , we obtain the encoded tokens  $\mathbf{x} = [x_0, x_1, \dots, x_i, \dots]$ , the perplexity of  $i^{th}$  token  $x_i$  in text  $S$  given model

<sup>3</sup>This name is inspired by the fact the sniffer dogs are able to trace the scents that cannot be easily noticed by humans.

$\theta_n$  is the log-likelihood  $ll_{\theta_n}(x_i) = \log p_{\theta_n}(x_i | x_{<i})$ . Given a list of known models  $\theta_0, \theta_1, \dots, \theta_N$ , we obtain a list of perplexities of the same text  $S$ . Since the tokenization process of each LLM might be different, we use a general word-level tokenization of  $S$ :  $\mathbf{w} = [w_0, w_1, \dots]$  and align calculated perplexities of tokens in  $\mathbf{x}$  to the general words. If the word  $w$  is aligned to multiple tokens in  $\mathbf{x}$ , we use the averaged perplexity; if a token in  $\mathbf{x}$  is aligned to multiple words, we assign these words with the same value.<sup>4</sup> Further, the aligned perplexity is conditioned on the specific model-training process, therefore, the perplexity should be normalized for comparison between models. We apply normalization strategies including dataset-wise-normalization and L1-normalization for these aligned features. Dataset-wise normalization is to normalize the perplexity with an averaged perplexity on all data. Therefore, as seen in Figure 2, we obtain lists of word-level perplexities  $ll_{\theta_n}(\mathbf{w})$  that are aligned across models.

##### Extract Contrastive Features

After obtaining aligned perplexities  $ll_{\theta_n}(\mathbf{w})$ , which is a list of token-wise perplexities, we aim to search features within the list of token-wise perplexities. Instead of introducing assumptions of the difference between human-written texts and machine-generated texts (Mitchell et al., 2023), the core idea of Sniffer is to find different features between models. Therefore, we pose a simple hypothesis:

**Hypothesis 1** *A human written text tends to have a similar perplexity list across models, and a generated text tends to show discrepancy across models.*

To further generalize the extracted features for tracing unknown models, we pose another hypothesis:

**Hypothesis 2** *The discrepancy of perplexity*

<sup>4</sup>The details of the aligning process can be seen in the Appendix.

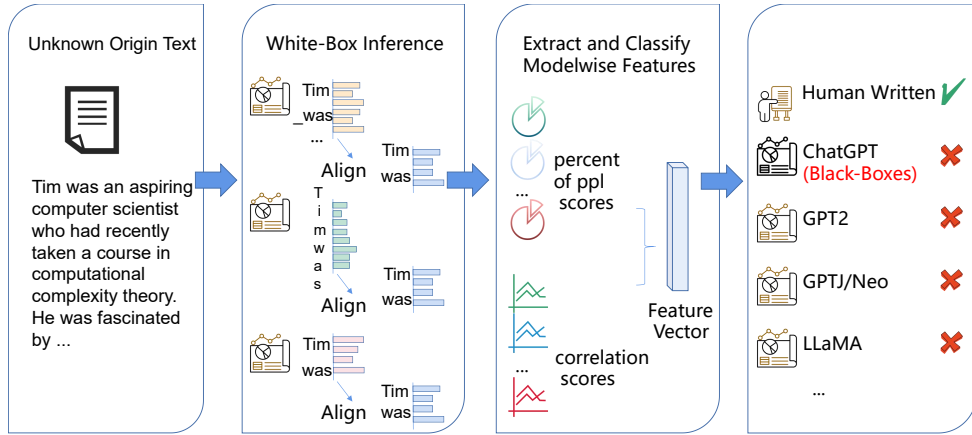


Figure 2: Sniffer process Illustration

257 *curves across models can reveal features of both*  
 258 *known and unknown models.*

259 With these hypotheses, we can utilize the  
 260 contrastive features between known models  
 261  $\theta_0, \theta_1, \dots$  to trace the origin of various models  
 262 denoted as  $\phi_0, \phi_1, \dots$  that includes known and  
 263 unknown models. Specifically, given a pair-wise  
 264 known model comparison between model  $\theta_i$   
 265 and  $\theta_j$ , we have  $ll_{\theta_i}(w)$  and  $ll_{\theta_j}(w)$  with  
 266 aligned perplexities of  $L$  words in total. We  
 267 first design a *percent-of-low-perplexity* score  
 268 that calculates the percentage of the lower word  
 269 perplexities of the perplexity of  $\theta_i$  compared with  
 270  $\theta_j$ :  $\text{pct}_{ij} = \left( \sum_k^N \mathbb{1}(ll_{\theta_i}(w_k) < ll_{\theta_j}(w_k)) \right) / L$ .  
 271 Here,  $\mathbb{1}(\cdot)$  is the indicator function that assigns 1  
 272 when the term is satisfied.

273 Given  $N$  known models, we have  $C_N^2$  pairs and  
 274 we can obtain  $C_N^2$   $\text{pct}_{ij}$  scores. We collect all  
 275 these scores as heuristic features that can be used  
 276 to trace the text origins. Plus, we also include  
 277 the sentence-level perplexities and the Pearson and  
 278 Spearman correlation coefficient between  $ll_{\theta_i}(w)$   
 279 and  $ll_{\theta_j}(w)$  to construct a feature vector as the  
 280 final representation for origin tracing. For instance,  
 281 given 4 known LLMs ( $N = 4$ ), we have 4 sentence-  
 282 level perplexity features,  $C_4^2 = 6$  pct-scores, and  
 283  $6 * 2 = 12$  correlation coefficient score, therefore  
 284 the final representation vector is a 22-dim vector.

### 285 Train Features for Origin Tracing

286 After extracting the contrastive features from  
 287 different known models, we can train a simple  
 288 linear classifier  $F(\cdot)$  to project the extracted  
 289 features to different model origins. To train the  
 290 classifier  $F(\cdot)$ , we collect a small amount of  
 291 human-written texts that cover various aspects and  
 292 use different known models to generate texts as  
 293 similar or *parallel* data compared to human-written  
 294 texts to train the linear classifier.

The trained classifier has several unique features:

(1) Low-resource required: Unlike semantic-  
 wise classification tasks that require a large amount  
 of data and a strong natural language encoder,  
 after feature extraction, the final representation is a  
 low-dimension vector, which requires only a small  
 amount of data since the feature contains abundant  
 model-wise heuristic knowledge for origin tracing;

(2) Generalization ability to unknown models  
 and stronger LLMs: In the classifier studying  
 process, we can collect unknown-model-generated  
 texts to obtain their features based on known  
 models and these features can reveal different  
 traces compared with known models and human-  
 written texts. Different from supervised-learning  
 methods that rely on semantic-level features to  
 detect text origins, model-wise features are  
**NOT** influenced as the quality of generated texts  
 improves. That is, in the supervised learning  
 methods, stronger LLMs are harder to detect since  
 the generated texts are human-like, but model-wise  
 features cannot be easily optimized.

(3) Extend Ability: In our proposed method,  
 the linear classifier can be easily modified to trace  
 various model origins. Given a new open-source  
 model, we can easily use the feature extraction  
 strategy above to re-train the classification model;  
 given a new unknown model, we can easily collect  
 a few generated texts and extract the features based  
 on known models, and train a new classifier with  
 the unknown model origin features to trace the new  
 unknown origin. With such extended ability, our  
 proposed method can be used in both black-box  
 and white-box settings.

329  
330  
331  
332  
333  
334  
335  
336  
337  
338  
339  
340  
341  
342  
343  
344  
345  
346  
347  
348  
349  
350  
351  
352  
353  
354  
355  
356  
357  
358  
359  
360  
361  
362  
363  
364  
365  
366  
367  
368  
369  
370  
371  
372  
373  
374  
375  
376  
377  
378  
379

## 4 Experiments

### 4.1 Dataset Construction

In the origin tracing of LLMs, one major challenge is that LLMs are almost omniscient to world knowledge since they are trained with various and huge amounts of data. We collect a wide range of texts from different origins for the proposed origin tracing tool Sniffer to serve as a general detector.

We collect texts from domains including News articles, social media posts, web texts, scientific articles or academic papers, and technical documentation. We use public datasets including XSum dataset (Narayan et al., 2018) that contains news articles; IMDB dataset that contains social media reviews of movies; web texts (Radford et al., 2019) that contains texts from common-crawled online pages; PubMed and Arxiv dataset (Cohan et al., 2018) that contains academic topics; Wikipedia corpus used in SQuAD dataset (Rajpurkar et al., 2016) that contains general world knowledge. For each dataset, we randomly collect 1,000 documents and shuffle them into a human-written text dataset containing 6k documents. Then we generate AI-generated contexts from the texts we collected as *parallel* data to study origin tracing.

When generating texts from language models such as GPT-2, we use the first 10 words as the prompt to generate a document. When generating texts from instruction-tuned models such as GPT-3.5(text-davinci-003) and ChatGPT (turbo), we give several instructions including re-write instruction and story generation instruction to obtain a similar AI-generated text. (We show instruction details in the Appendix.) With these simple instructions, we collect AI-generated texts from instruction-tuned LLMs. For known origin models, we collect AI-generated texts from GPT2 (powered by OpenAI), GPT-J and GPT-Neo (powered by EleutherAI), LLaMA (powered by Meta AI). For unknown origin models, we collect AI-generated texts from GPT3.5-text-davinci-003, which is an instruction-tuned model with 175B parameters. Therefore, we collect 6k human written texts, and every 6k texts from GPT2, GPT-Neo, GPT-J, and 12k texts from GPT3 models, which is 36k texts in total. We divide the 36k texts into a train/test split with a 90%/10% partition.

We name the collected dataset **SnifferBench**, which can be further used in origin tracing and

AI-generated contexts detection tasks. Further, the dataset collection process can be extended to different LLMs in the future with more different scenarios and prompts/instructions which further challenges the origin tracing ability.

### 4.2 Implementations of Sniffer

In Sniffer, we select several open-source (L)LMs as known models: we use GPT2-xl(1.5B), GPT-Neo(2.7B), GPT-J(6B) and LLaMA(7B) as known models. In the SnifferBench, we collect texts from origins including GPT2(OpenAI), GPT-Neo and GPT-J(EleutherAI), LLaMA (MetaAI), ChatGPT(GPT3.5-turbo from OpenAI), and human-written texts, therefore, the unknown origin is the ChatGPT(GPT3.5-turbo) model since they are not open-source models. The goal of origin tracing in Sniffer is to trace both known and unknown origins. We construct an inference server for each known model based on NVIDIA4090 GPUs and set the max sequence length given the maximum GPU allowance. We align all texts with a white-space tokenizer to obtain uniform tokenizations.

We utilize the aforementioned four known models and employ a linear classifier to construct a set of binary and multi-class Sniffers, used respectively for comparison with existing binary models like DetectGPT, and for origin tracing. In addition, we introduce a new variant of Sniffer: **Sniffer(+GPT3)**. As we are able to obtain logits of the generated texts in the GPT-3 API provided by OpenAI, we are able to treat GPT-3 model as a white-box model. Therefore, we collect a subset in the SnifferBench to test the origin tracing with GPT-3 models as white boxes. Therefore, in our implementation of Sniffer, we have 5 models in total which result in a 35-dimension vector in obtaining the sniffer feature. Here, we align the tokens with the tokenizer provided by OpenAI.

### 4.3 Metrics of Origin Tracing

In the SnifferBench, we collect texts from different origins, therefore, we calculate the precision and recall of each text origin tracing prediction. We can only calculate the known model and the human-written texts' precision and recall in baseline methods since these methods cannot be used in the origin tracing task.

380  
381  
382  
383  
384  
385  
386  
387  
388  
389  
390  
391  
392  
393  
394  
395  
396  
397  
398  
399  
400  
401  
402  
403  
404  
405  
406  
407  
408  
409  
410  
411  
412  
413  
414  
415  
416  
417  
418  
419  
420  
421  
422  
423  
424  
425  
426

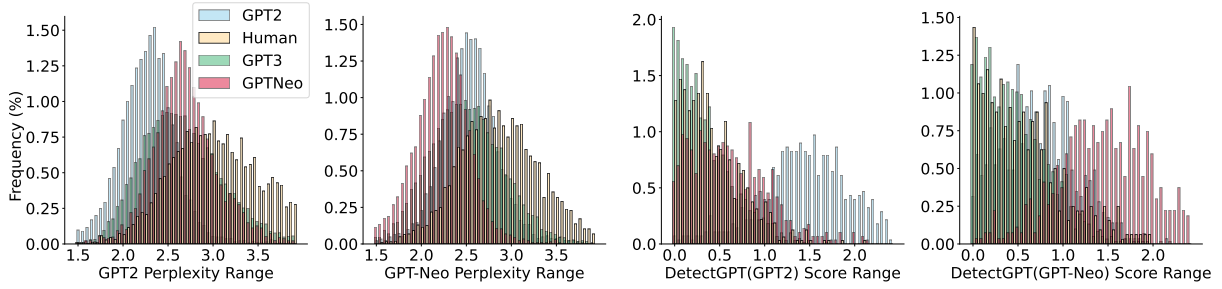


Figure 3: The discrepancy between different text origins in different baseline methods. In each figure, different bars show different text origins and each figure is to use a certain model of a certain detect method to test given texts.

Method	Different Text Origins					
	AI	GPT-2 Human	Overall	AI	GPT-Neo Human	Overall
$\log p(x)$	80.9/89.3	87.9/78.7	84.0	71.7/78.9	76.3/68.5	73.6
DetectGPT	88.9/88.9	89.9/90.2	89.5	74.4/79.3	80.0/75.5	77.2
Sniffer	<b>99.1/99.3</b>	<b>99.3/99.2</b>	<b>99.2</b>	<b>99.7/99.0</b>	<b>99.5/99.8</b>	<b>99.4</b>

Table 2: Precision/Recall Results of binary detection, we use macro-F1 as the overall metric. Since  $\log p(x)$  and DetectGPT methods are unable to trace different text origins, only the corresponding model-generated text and human-written text detection results are listed.

#### 4.4 Results of Binary Detection

As we compare Sniffer with previous detection methods in the form of the methods, we construct experiments to explore how previous methods fail to run origin tracing.

We first implement  $\log p(x)$  method which is also used in the original GPTZero<sup>5</sup>; then we implement DetectGPT method which introduces a discrepancy score to discriminate AI-generated texts by adding multiple perturbations (we try 40 perturbations for each sample). We collect the sentence-level perplexity and the DetectGPT discrepancy score of the collected dataset based on a certain known model  $\theta_n$  and draw a histogram showing the score distributions of different origin texts.

In Figure 3, we plot the histogram of the discrepancy between different texts’ origins in  $\log p(x)$  and DetectGPT. In  $\log p(x)$ , we plot the perplexity score of different text origins using a specific model such as GPT-2 or GPT-Neo, and in the DetectGPT method, we use the z-score proposed in DetectGPT. As seen in Figure 3, though there are multiple peaks showing that there are differences between different text origins, the overlap is too large to successfully separate different text origins. we can conclude that it is extremely difficult to discriminate text origins from

<sup>5</sup><https://gptzero.me/>

$\log p(x)$  features or z-score features. Therefore, it is important to introduce strong features to trace the texts’ origins.

Since  $\log p(x)$  and DetectGPT methods are unable to trace different text origins, we only conduct the corresponding model-generated text and human-written text detection. Specifically, we select a threshold manually as the discrimination boundary in  $\log p(x)$  and DetectGPT. As listed in Table 2, even in the binary classification detection tasks they excel at, their performance still falls short of Sniffer. This further confirms that whether in origin tracing capability or binary classification detection ability, Sniffer is an exceptional method.

#### 4.5 Results of Sniffer Origin Tracing

In Table 3, we can observe that Sniffer can obtain impressive results in tracing the text origins in the known models including GPT-2, GPT-J/Neo; further, when Sniffer does not use the GPT-3 models as white boxes, Sniffer can generalize to trace GPT-3 model origins and can properly discriminate them from human-written lines. Such an ability cannot be easily obtained as the GPT-3 models are stronger and similar to humans, meanwhile remaining unknown to the public.

Plus, we can observe that recently released LLaMA models are more difficult to trace, though it is a white box in Sniffer, indicating that stronger LLMs are harder to detect. When LLMs are being

Method	Different Text Origins					Overall
	GPT-2 (OpenAI)	GPT-J/Neo (EleutherAI)	LLama (Meta)	ChatGPT (OpenAI)	Human	
Random	16.6/16.1	16.6/16.1	16.6/16.1	16.6/16.1	16.6/16.1	16.3
Sniffer	98.7/96.9	96.6/ <b>98.0</b>	<b>85.0/84.3</b>	<b>77.7/82.3</b>	68.1/60.3	<b>84.8</b>
Sniffer (10%)	97.3/96.3	<b>96.7/96.1</b>	80.9/77.2	73.9/77.3	58.9/ <b>67.7</b>	82.2
Sniffer (5%)	97.3/97.5	96.6/95.1	76.1/74.0	71.4/76.7	58.8/53.4	79.7
Sniffer (1%)	97.9/94.4	91.0/95.2	65.8/60.2	67.4/76.3	60.0/46.4	75.4
Sniffer (L1-norm)	<b>97.8/98.3</b>	<b>96.7/95.9</b>	75.2/74.4	74.7/82.4	<b>75.7/62.4</b>	83.3
Sniffer ( $\log p(x)$ only)	<b>98.9/97.7</b>	94.1/94.8	60.4/49.3	64.6/78.8	63.0/47.6	74.9
Sniffer (pct-score only)	98.3/96.6	94.0/94.8	59.5/53.2	60.0/79.9	58.3/26.8	72.1
Sniffer ( $\log p(x)$ + pct-score)	98.6/97.2	96.5/96.2	69.6/65.0	71.0/ <b>82.5</b>	66.3/51.1	79.4

Table 3: Precision/Recall Results of Origin Tracing, we use macro-F1 as the overall metric. % in Sniffer is to train Sniffer with limited training data;  $\log p(x)$  and pct-score -only are the variants that only use the corresponding features for the linear projection.

Method	Different Text Origins				Human	Overall
	GPT-2 (OpenAI)	GPT-J/Neo (EleutherAI)	LLama (Meta)	GPT-3 (OpenAI)		
Sniffer(1%)	<b>93.8/90.0</b>	89.6/88.0	59.1/47.4	57.6/59.6	35.3/42.9	66.3
Sniffer(+GPT3)(1%)	92.0/ <b>91.8</b>	<b>90.0/92.4</b>	<b>81.6/72.6</b>	<b>70.0/65.6</b>	<b>42.1/50.0</b>	<b>74.8</b>

Table 4: Results of Sniffer(+GPT3): We use GPT-3.5(text-davinci-003) to replace ChatGPT(turbo) since ChatGPT(turbo) does not provide logits.

studied by a wide range of researchers, it is also important to be alert that stronger models may be harder to detect, and harder to control which can cause potential harm to society.

## 4.6 Different Sniffers

### Sniffer+GPT3

In Table 4, we list the results of using Sniffer(+GPT3) to trace the origins of texts generated by GPT3 models while we treat GPT3 models as white boxes. We use the OpenAI service that returns token logits therefore we use the text-davinci-003 model and test on a 1% dataset compared with the full dataset. Here, we replace texts generated from ChatGPT with davinci-003 outputs therefore the testset is different from the one used in tracing ChatGPT texts.

As seen, when the Sniffer feature can use GPT-3 logits, the results grow significantly higher compared with generalizing features from GPT2/J/Neo and LLaMA models to trace GPT-3. Therefore, we can conclude that although Sniffer is able to generalize its features to trace unknown origins, it is better to possess the LLMs, indicating that the actual LLM providers must be more cautious when releasing LLMs since they are more capable to avoid abuse or malicious usage of LLMs.

### Limited Data Experiments

As we illustrated, in methods using model-wise features, we can use limited data to construct a powerful detector. Therefore, we use different numbers of training data to train Sniffer. As seen, when we use limited data to train the Sniffer model, the performances are not significantly harmed, indicating that the model-wise features are high-quality features that reveal obvious traces of texts for the origin tracing classification.

### Feature Ablations

In the extracted features, we observe that using L1-norm can help obtain a higher performance in tracing human texts, but is rather weak in tracing different model origins, especially when tracing LLaMA. Therefore, we use datase-wise norm for the rest experiments.

To further analyze how the extracted features help trace the text origins, we run a simple ablation test that uses different features proposed in Sniffer. As seen in Table 3, the perplexity score  $\log p(x)$  is one key metric but can be significantly improved by the percent-of-perplexity score (pct-score), showing that we can trace text origins by analyzing the discrepancies between models when testing same texts as Sniffer did.

## 4.7 Difficulty of Tracing Different Types of Generated Texts

### Generation Genre Influence

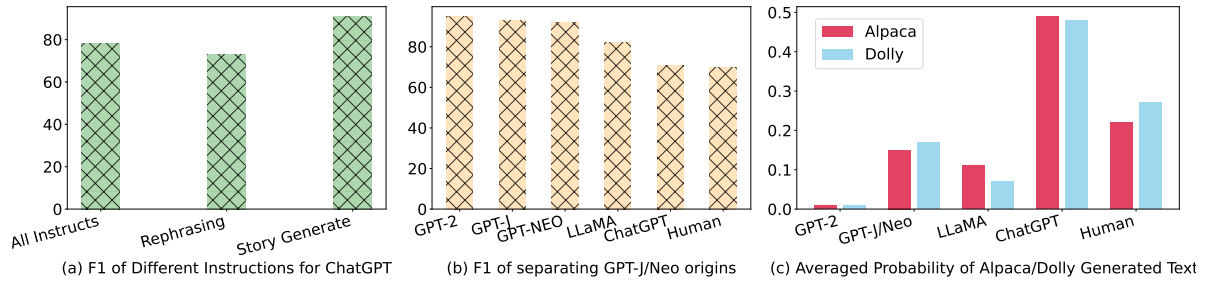


Figure 4: Tracing different types of generated texts: (a) plots the ChatGPT tracing results that use different instructions; (b) plots the tracing results that separate GPT-J and GPT-Neo origins and calculate the f1-score of corresponding origins; (c) plots the tracing results that test Alpaca and Dolly models that use ChatGPT instructions to supervise fine-tuning the LLaMA/GPT-J models to build an instructed LLM.

As illustrated, we use several different instructions to instruct GPT3.5 (turbo) models and GPT3.5-text-davinci-003 models, we further discuss the tracing difficulty when the texts are generated by different instructions.

As seen in Figure 4(a), we find that when the texts are generated by rephrasing instructions, the texts are rather easy to trace while texts generated from a summary are harder to trace. Such an observation indicates that the difficulty of tracing texts generated by strong LLMs is also different when the instructions are different. The rephrased texts are more similar to human-written texts, therefore, are more difficult to detect.

### Mixing Origin Tracing

In our origin tracing experiment setup, we divide GPT-J/Neo generated texts into the same origin since these models are provided by the same facility. That is, LLM origins can have different levels of classifications. As LLMs can be trained with *distill texts* from strong LLMs such as ChatGPT, it is harder to trace the origin if the model is trained based on a base model such as GPT-J or LLaMA but is instructed by ChatGPT-generated outputs. It is hard to tell whether the base model or the instruct model has a larger impact on the mixed-origin model. Therefore, we construct experiments to first separate GPT-J/Neo origins using the proposed datasets. Further, we test mixed-origin models such as Alpaca (Taori et al., 2023) (ChatGPT instructions tuned based on LLaMA) and Dolly<sup>6</sup> (ChatGPT instructions tuned based on GPT-J) by instructing these models to generate 400 samples and testing them with Sniffer.

In Figure 4(b), we show the f1-score of origin

tracing that separates GPT-J and GPT-Neo models using the full data of SnifferBench. As seen, Sniffer is able to successfully divide GPT-J and GPT-Neo, indicating that the text origin divide can be of different levels. We can trace the text origins of a specific model or a certain party.

In Figure 4(c), we list the averaged probability of the Sniffer inference results of Alpaca and Dolly. As seen, texts from both Alpaca and Dolly tend to be categorized as ChatGPT-generated texts, indicating that the align process has a more significant impact on the generated texts compared with the base model. Therefore, Sniffer can be used as a detector for testing whether a model is trained from ChatGPT-generated instructions, helping protect the originality of LLMs.

## 5 Conclusion and Future Work

In this paper, we first introduce the concept of origin tracing, an important direction in the era of LLMs. Then we discuss two lines of AI-generated context detection and origin tracing methods and point out the necessity of studying model-wise features for origin tracing. We further design a simple Sniffer method as well as a benchmark to test the origin tracing challenge. Through extensive experiments, we find that the current origin tracing field is full of challenges including tracing texts from models with mixed origins; combining semantic features and model-wise features; and tracing texts from LLMs that are given various instructions. Therefore, we can hope for a continuous line of work that studies the origin tracing of LLMs and hopes to improve the trustworthy and safe usage of LLMs.

<sup>6</sup><https://huggingface.co/databricks/dolly-v2-12b>



## 609 Limitations

610 Despite Sniffer having achieved substantial origin  
611 tracing results and unearthing many valuable  
612 insights, there are still some limitations:

613 (1) This study exclusively employs a simple  
614 linear classifier for origin tracing, without the  
615 utilization of more intricate models. This approach  
616 is adopted as we pioneer the application of model-  
617 wise perplexity features for origin tracing and  
618 AI-generated text detection, and it has proven  
619 sufficiently effective even on a linear classifier.  
620 Future research will explore the use of more  
621 complex models to enhance performance further.

622 (2) With more complicated model-generated  
623 texts, classifiers based on semantics can achieve  
624 better results. In the future, we plan to integrate  
625 Sniffer with semantic-wise models to experiment  
626 with more robust models.

627 As mentioned, we align tokens with different  
628 tokenization methods with a uniform tokenization  
629 strategy to compare the perplexity at the same  
630 level. As shown in Figure 6 in the Appendix,  
631 when the texts are tokenized by different tokenizers,  
632 we project different tokens to the byte level, and  
633 from the byte-level projections, we align different  
634 tokenized texts to a list of tokens with the same  
635 tokenizations.

## 636 References

637 Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda  
638 Askell, Anna Chen, Nova DasSarma, Dawn Drain,  
639 Stanislav Fort, Deep Ganguli, T. J. Henighan,  
640 Nicholas Joseph, Saurav Kadavath, John Kernion,  
641 Tom Conerly, Sheer El-Showk, Nelson Elhage, Zac  
642 Hatfield-Dodds, Danny Hernandez, Tristan Hume,  
643 Scott Johnston, Shauna Kravec, Liane Lovitt, Neel  
644 Nanda, Catherine Olsson, Dario Amodei, Tom B.  
645 Brown, Jack Clark, Sam McCandlish, Christopher  
646 Olah, Benjamin Mann, and Jared Kaplan. 2022.  
647 Training a helpful and harmless assistant with  
648 reinforcement learning from human feedback. *ArXiv*,  
649 abs/2204.05862.

650 Anton Bakhtin, Sam Gross, Myle Ott, Yuntian Deng,  
651 Marc’Aurelio Ranzato, and Arthur D. Szlam. 2019.  
652 Real or fake? learning to discriminate machine from  
653 human generated text. *ArXiv*, abs/1906.03351.

654 Sid Black, Leo Gao, Phil Wang, Connor Leahy, and  
655 Stella Rose Biderman. 2021. Gpt-neo: Large  
656 scale autoregressive language modeling with mesh-  
657 tensorflow.

658 Tom Brown, Benjamin Mann, Nick Ryder, Melanie  
659 Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind

Neelakantan, Pranav Shyam, Girish Sastry, Amanda  
Askell, et al. 2020. Language models are few-shot  
learners. *Advances in neural information processing  
systems*, 33:1877–1901. 660 661 662 663

Nicholas Carlini, Florian Tramèr, Eric Wallace,  
Matthew Jagielski, Ariel Herbert-Voss, Katherine  
Lee, Adam Roberts, Tom B. Brown, Dawn Xiaodong  
Song, Úlfar Erlingsson, Alina Oprea, and Colin  
Raffel. 2020. Extracting training data from large  
language models. In *USENIX Security Symposium*. 664 665 666 667 668 669

Arman Cohan, Franck Dernoncourt, Doo Soon  
Kim, Trung Bui, Seokhwan Kim, Walter Chang,  
and Nazli Goharian. 2018. [A discourse-aware  
attention model for abstractive summarization of  
long documents](#). In *Proceedings of the 2018  
Conference of the North American Chapter of the  
Association for Computational Linguistics: Human  
Language Technologies, Volume 2 (Short Papers)*,  
pages 615–621, New Orleans, Louisiana. Association  
for Computational Linguistics. 670 671 672 673 674 675 676 677 678 679

Yiming Cui, Ting Liu, Wanxiang Che, Li Xiao, Zhipeng  
Chen, Wentao Ma, Shijin Wang, and Guoping Hu.  
2019. [A span-extraction dataset for Chinese machine  
reading comprehension](#). In *Proceedings of the  
2019 Conference on Empirical Methods in Natural  
Language Processing and the 9th International  
Joint Conference on Natural Language Processing  
(EMNLP-IJCNLP)*, pages 5883–5889, Hong Kong,  
China. Association for Computational Linguistics. 680 681 682 683 684 685 686 687 688

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and  
Kristina Toutanova. 2018. [BERT: pre-training  
of deep bidirectional transformers for language  
understanding](#). *CoRR*, abs/1810.04805. 689 690 691 692

Brian Dolhansky, Joanna Bitton, Ben Pflaum, Jikuo Lu,  
Russ Howes, Menglin Wang, and Cristian Canton-  
Ferrer. 2020. The deepfake detection challenge  
dataset. *ArXiv*, abs/2006.07397. 693 694 695 696

Sebastian Gehrmann, Hendrik Strobelt, and Alexander  
Rush. 2019. [GLTR: Statistical detection and  
visualization of generated text](#). In *Proceedings  
of the 57th Annual Meeting of the Association for  
Computational Linguistics: System Demonstrations*,  
pages 111–116, Florence, Italy. Association for  
Computational Linguistics. 697 698 699 700 701 702 703

Jian Guan, Zhuoer Feng, Yamei Chen, Ruilin He, Xiaoxi  
Mao, Changjie Fan, and Minlie Huang. 2022. Lot:  
A story-centric benchmark for evaluating chinese  
long text understanding and generation. *Transactions  
of the Association for Computational Linguistics*,  
10:434–451. 704 705 706 707 708 709

Ganesh Jawahar, Muhammad Abdul-Mageed, and Laks  
V. S. Lakshmanan. 2020. Automatic detection  
of machine generated text: A critical survey.  
In *International Conference on Computational  
Linguistics*. 710 711 712 713 714

715	John Kirchenbauer, Jonas Geiping, Yuxin Wen,	<a href="#">transformer</a> . <i>Journal of Machine Learning Research</i> ,	770
716	Jonathan Katz, Ian Miers, and Tom Goldstein. 2023.	21(140):1–67.	771
717	A watermark for large language models. <i>ArXiv</i> ,		
718	abs/2301.10226.		
719	Keita Kurita, Paul Michel, and Graham Neubig. 2020.	Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and	772
720	Weight poisoning attacks on pretrained models.	Percy Liang. 2016. Squad: 100,000+ questions	773
721	<i>ArXiv</i> , abs/2004.06660.	for machine comprehension of text. <i>arXiv preprint</i>	774
		<i>arXiv:1606.05250</i> .	775
722	Linyang Li, Ruotian Ma, Qipeng Guo, Xiangyang Xue,	Irene Solaiman, Miles Brundage, Jack Clark, Amanda	776
723	and Xipeng Qiu. 2020. Bert-attack: Adversarial	Askell, Ariel Herbert-Voss, Jeff Wu, Alec Radford,	777
724	attack against bert using bert. <i>arXiv preprint</i>	and Jasmine Wang. 2019. Release strategies and	778
725	<i>arXiv:2004.09984</i> .	the social impacts of language models. <i>ArXiv</i> ,	779
		abs/1908.09203.	780
726	Linyang Li, Demin Song, Xiaonan Li, Jiehang Zeng,	Chakraborty Souradip, Singh Bedi, Amrit, Zhu Sicheng,	781
727	Ruotian Ma, and Xipeng Qiu. 2021. Backdoor	An Bang, Manocha, Dinesh, and Huang Furong.	782
728	attacks on pre-trained models by layerwise weight	2023. On the possibilities of ai-generated text	783
729	poisoning. In <i>Conference on Empirical Methods in</i>	detection.	784
730	<i>Natural Language Processing</i> .		
731	Yudong Li, Yuqing Zhang, Zhe Zhao, Linlin	Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann	785
732	Shen, Weijie Liu, Weiquan Mao, and Hui	Dubois, Xuechen Li, Carlos Guestrin, Percy Liang,	786
733	Zhang. 2022. <a href="#">CSL: A large-scale Chinese</a>	and Tatsunori B. Hashimoto. 2023. Stanford alpaca:	787
734	<a href="#">scientific literature dataset</a> . In <i>Proceedings of the</i>	An instruction-following llama model. <a href="https://github.com/tatsu-lab/stanford_alpaca">https://</a>	788
735	<i>29th International Conference on Computational</i>	<a href="https://github.com/tatsu-lab/stanford_alpaca">github.com/tatsu-lab/stanford_alpaca</a> .	789
736	<i>Linguistics</i> , pages 3917–3923, Gyeongju, Republic		
737	of Korea. International Committee on Computational	Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier	790
738	Linguistics.	Martinet, Marie-Anne Lachaux, Timothée Lacroix,	791
		Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal	792
739	Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du,	Azhar, Aurélien Rodriguez, Armand Joulin, Edouard	793
740	Mandar Joshi, Danqi Chen, Omer Levy, Mike	Grave, and Guillaume Lample. 2023. Llama: Open	794
741	Lewis, Luke Zettlemoyer, and Veselin Stoyanov.	and efficient foundation language models. <i>ArXiv</i> ,	795
742	2019. Roberta: A robustly optimized bert pretraining	abs/2302.13971.	796
743	approach. <i>arXiv preprint arXiv:1907.11692</i> .		
744	Eric Mitchell, Yoonho Lee, Alexander Khazatsky,	Ben Wang. 2021. Mesh-Transformer-JAX: Model-	797
745	Christopher D. Manning, and Chelsea Finn.	Parallel Implementation of Transformer Language	798
746	2023. Detectgpt: Zero-shot machine-generated	Model with JAX. <a href="https://github.com/kingoflolz/mesh-transformer-jax">https://github.com/</a>	799
747	text detection using probability curvature. <i>ArXiv</i> ,	<a href="https://github.com/kingoflolz/mesh-transformer-jax">kingoflolz/mesh-transformer-jax</a> .	800
748	abs/2301.11305.		
749	Shashi Narayan, Shay B. Cohen, and Mirella Lapata.	Liang Xu, Hai Hu, Xuanwei Zhang, Lu Li, Chenjie Cao,	801
750	2018. Don’t give me the details, just the summary!	Yudong Li, Yechen Xu, Kai Sun, Dian Yu, Cong	802
751	Topic-aware convolutional neural networks for	Yu, Yin Tian, Qianqian Dong, Weitang Liu, Bo Shi,	803
752	extreme summarization. In <i>Proceedings of the</i>	Yiming Cui, Junyi Li, Jun Zeng, Rongzhao Wang,	804
753	<i>2018 Conference on Empirical Methods in Natural</i>	Weijian Xie, Yanting Li, Yina Patterson, Zuoyu Tian,	805
754	<i>Language Processing</i> , Brussels, Belgium.	Yiwen Zhang, He Zhou, Shaowei Hua Liu, Zhe Zhao,	806
		Qipeng Zhao, Cong Yue, Xinrui Zhang, Zhengliang	807
755	OpenAI. 2023. Gpt-4 technical report. <i>ArXiv</i> ,	Yang, Kyle Richardson, and Zhenzhong Lan. 2020.	808
756	abs/2303.08774.	<a href="#">CLUE: A Chinese language understanding evaluation</a>	809
		<a href="#">benchmark</a> . In <i>Proceedings of the 28th International</i>	810
757	Alec Radford, Karthik Narasimhan, Tim Salimans,	<i>Conference on Computational Linguistics</i> , pages	811
758	and Ilya Sutskever. 2018. Improving language	4762–4772, Barcelona, Spain (Online). International	812
759	understanding by generative pre-training. <i>URL</i>	Committee on Computational Linguistics.	813
760	<a href="https://s3-us-west-2.amazonaws.com/openai-assets/researchcovers/languageunsupervised/language">https://s3-us-west-2.</a>	Rowan Zellers, Ari Holtzman, Hannah Rashkin,	814
761	<a href="https://s3-us-west-2.amazonaws.com/openai-assets/researchcovers/languageunsupervised/language">amazonaws.</a>	Yonatan Bisk, Ali Farhadi, Franziska Roesner, and	815
762	<a href="https://s3-us-west-2.amazonaws.com/openai-assets/researchcovers/languageunsupervised/language">com/openai-</a>	Yejin Choi. 2019. Defending against neural fake	816
	<a href="https://s3-us-west-2.amazonaws.com/openai-assets/researchcovers/languageunsupervised/language">assets/researchcovers/languageunsupervised/language</a>	news. In <i>Neural Information Processing Systems</i> .	817
	<a href="https://s3-us-west-2.amazonaws.com/openai-assets/researchcovers/languageunsupervised/language">understanding paper.pdf</a> .		
763	Alec Radford, Jeff Wu, Rewon Child, David Luan,	Aohan Zeng, Xiao Liu, Zhengxiao Du, Zihan Wang,	818
764	Dario Amodei, and Ilya Sutskever. 2019. Language	Hanyu Lai, Ming Ding, Zhuoyi Yang, Yifan Xu,	819
765	models are unsupervised multitask learners.	Wendi Zheng, Xiao Xia, Weng Lam Tam, Zixuan Ma,	820
		Yufei Xue, Jidong Zhai, Wenguang Chen, Zhiyuan	821
766	Colin Raffel, Noam Shazeer, Adam Roberts, Katherine	Liu, Peng Zhang, Yuxiao Dong, and Jie Tang. 2023.	822
767	Lee, Sharan Narang, Michael Matena, Yanqi Zhou,	<a href="#">GLM-130b: An open bilingual pre-trained model</a> . In	823
768	Wei Li, and Peter J. Liu. 2020. <a href="#">Exploring the</a>	<i>The Eleventh International Conference on Learning</i>	824
769	<a href="#">limits of transfer learning with a unified text-to-text</a>	<i>Representations (ICLR)</i> .	825

826 Jiaxing Zhang, Ruyi Gan, Junjie Wang, Yuxiang Zhang,  
827 Lin Zhang, Ping Yang, Xinyu Gao, Ziwei Wu,  
828 Xiaoqun Dong, Junqing He, Jianheng Zhuo, Qi Yang,  
829 Yongfeng Huang, Xiayu Li, Yanghan Wu, Junyu Lu,  
830 Xinyu Zhu, Weifeng Chen, Ting Han, Kunhao Pan,  
831 Rui Wang, Hao Wang, Xiaojun Wu, Zhongshen Zeng,  
832 and Chongpei Chen. 2022. Fengshenbang 1.0: Being  
833 the foundation of chinese cognitive intelligence.  
834 *CoRR*, abs/2209.02970.

835 Jieyu Zhao, Tianlu Wang, Mark Yatskar, Vicente  
836 Ordonez, and Kai-Wei Chang. 2017. Men also  
837 like shopping: Reducing gender bias amplification  
838 using corpus-level constraints. *arXiv preprint*  
839 *arXiv:1707.09457*.

Table 5: Precision/Recall results of Chinese Dataset Performances of Origin Tracing.

Method	Different Text Origins						Human	Overall
	Wenzhong Fengshenbang	Damo Damo	Skytext SkyWork	ChatGLM THUDM	ChatGPT OpenAI	MOSS FudanNLP		
Sniffer	85.5/89.0	89.4/89.8	95.2/95.5	81.5/85.7	75.9/82.0	69.7/50.4	60.1/56.4	79.0

Table 6: Details of instructions used in Generating texts from ChatGPT, GPT3.5(text-davinci-003), Alpaca and Dolly models. In the story-generation, we use two instructs to obtain the generated document.

Instruction Type	Models	Instructions
Re-write	ChatGPT	You are an assistant that can rephrase the provided content. rephrase the following content: [document]
Story-Gen	ChatGPT	<i>Instruct 1:</i> You are an assistant that can summarize the provided content. summarize the following content: [document] <i>Instruct 2:</i> You are an assistant that can write a natural and fluent document based on the provided content, do not use fixed or academic writing style write a natural and fluent document like human based on the following content: [summary]
Re-write	GPT3.5	Rephrase the following content: [document]
Re-write	Alpaca & Dolly	Below is an instruction that describes a task, paired with an input that provides further context. Write a response that appropriately completes the request.  ### Instruction: Rewrite the following paragraph in a different style using your own words.  ### Input: [document]

## A Appendix

### Chinese Corpus Experiments

We also conducted a Chinese version of Sniffer as well as a collected testset for studying different types of LLMs.

We sample 1k samples from each open-source dataset of various domains including a long-document corpus of stories (Guan et al., 2022); various web texts from Iflytek classification dataset (Xu et al., 2020); Chinese academical documents from CSL dataset (Li et al., 2022); Chinese Wikipedia corpus from CMRC dataset (Cui et al., 2019); review corpus from Xiecheng APP<sup>7</sup>. We collect these datasets that cover various domains with relatively long documents (more than 200 Chinese characters per document). The dataset construction setup is similar to English dataset setups. We also use ChatGPT (GPT turbo) model to generate Chinese texts with the Chinese instructions as ChatGPT (GPT turbo) is a strong multi-lingual model.

<sup>7</sup><https://huggingface.co/datasets/seamew/ChnSentiCorp>

As for open-source model selections, we adopt several open-source Chinese LLMs including the Wenzhong model (Zhang et al., 2022), a 2.7B GPT-2 style Chinese LLM; Damo model<sup>8</sup>, another 2.7B GPT-2 style Chinese LLM; Skytext model<sup>9</sup>, a 3B chatbot, and ChatGLM (Zeng et al., 2023), a 6B ChatGPT-style model trained with instructions.

For the generalization test of black-box models, we use ChatGPT (GPT turbo) and MOSS, a 16B Chinese LLM<sup>10</sup> with instructions that ask LLMs to re-write the given document as black-box origin tests.

We generate Chinese AI-generated contexts with open-source models and black-box models to build the whole Chinese benchmark for origin tracing tests.

As seen in Table 5, the Chinese texts show similar performances with English texts, indicating that the origin tracing strategy can be used in

<sup>8</sup>[https://modelscope.cn/models/damo/nlp\\_gpt3\\_text-generation\\_2.7B/summary](https://modelscope.cn/models/damo/nlp_gpt3_text-generation_2.7B/summary)

<sup>9</sup><https://huggingface.co/SkyWork/SkyTextTiny>

<sup>10</sup><https://github.com/OpenLMLab/MOSS>

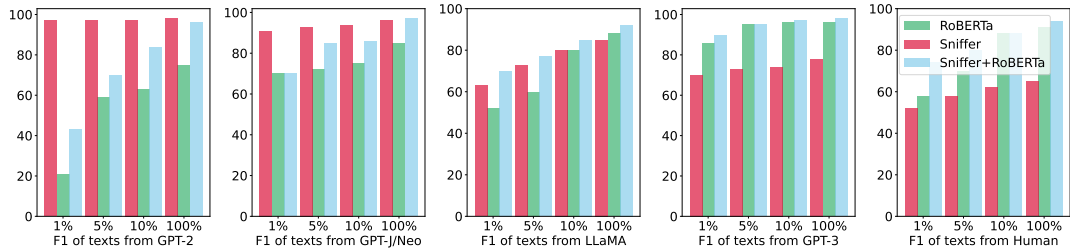


Figure 5: Comparison with Supervised Learning Methods that utilize semantic-wise features.

different languages.

### Instruction Details

In table 6, we list the specific instructions used in LLMs. We consider simple instructions since we focus on studying the origin tracing of LLMs without further analyzing the working mechanisms and better usage of instructions, which we leave to future works.

### Compare with Supervised Learning Methods

As discussed, the model-wise features trace the text origins by analyzing the model discrepancies, and the semantic-wise features are used as a fine-tuning task. We compare Sniffer with a fine-tuned RoBERTa model and measure the f1-score. Further, we combine Sniffer features and RoBERTa [CLS] feature to train a linear classifier as Sniffer+RoBERTa to test the performances of combinations of model-wise and semantic-wise features. As seen in Figure 5, the supervised learning method fails to solve the problem when the data is limited. Still, when the data is abundant, the performances are stronger than model-wise features (Mitchell et al., 2023; Souradip et al., 2023).

We need to notice that the supervised learning method fails to detect GPT-2/GPT-J/Neo generated texts, which is different from Sniffer results as GPT-2 texts are known to be less fluent compared with stronger models such as GPT-3 models. Therefore, we can assume that the RoBERTa models may capture some specific patterns since the GPT-3 generated texts only use several static instructions and the semantic features are easier to detect. The failure of GPT-2 detection indicates that semantic features are limited in origin tracing, calling for better algorithms to utilize model-wise features. In Sniffer+RoBERTa, we obtain promising results in tracing all origins, indicating that a proper combination of two types of features can help better trace text origins.

### Case Studies

We list several case studies that are randomly selected from the testset. We show the texts to be

detected, percent-of-perplexity scores, correlation scores calculated by Sniffer, and the tracing results.

As listed below, we show the list of perplexity scores of the given text and the extracted features, then we show the predicted tracing result. As seen, the perplexity list shows a similar trend between different models, but the perplexity value tends to show differences across models. For instance, in the GPT-2 generated text calculated list of perplexity, the GPT-2 perplexity is relatively lower than other perplexities calculated by other models, revealing the feature that helps trace the GPT-2 texts. On the other hand, for human-written texts, and texts generated by strong LLMs such as ChatGPT, the value does not show obvious differences between different models, for instance, some peaks in the perplexity list can be from the LLaMA models and some can be from some other models, which supports the hypothesis made above.

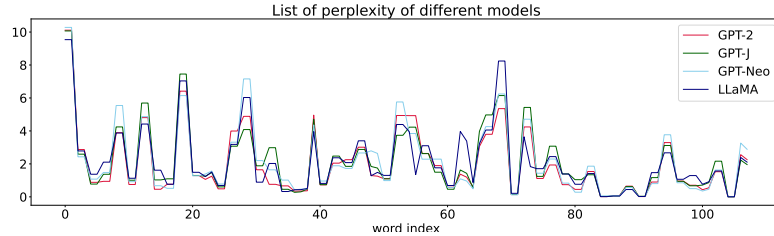
---

Text Origin

---

I thought that Mukhsin has been wonderfully written. Its not a work of fiction, but it certainly feels like a novel. I agree that there is much to admire about, despite flaws, which will be evident to many of its readers.

GPT-2



---

Extracted Features:

$\log p(x)$  (GPT-2): 2.04  
 $\log p(x)$  (GPT-Neo):2.18  
 $\log p(x)$  (GPT-J): 2.18  
 $\log p(x)$  (LLaMA): 2.21

pct-scores / Pearson scores / Spearman scores:

GPT-2 and GPT-Neo	0.59/0.97/0.97
GPT-2 and GPT-J	0.57/0.96/0.95
GPT-2 and LLaMA	0.63/0.92/0.91
GPT-Neo and GPT-J	0.56/0.94/0.94
GPT-Neo and LLaMA	0.52/0.92/0.89
GPT-J and LLaMA	0.55/0.92/0.93

---

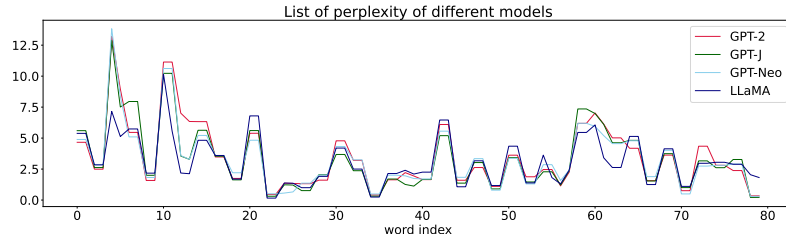
Tracing Results:

GPT-2: 99.8 %  
GPT-J/Neo: 0.1 %  
LLaMA: 0.0%  
ChatGPT: 0.1 %  
Human: 0.0 %

---

Ten passengers on board the bus were reported to have either been injured, or have since been released from hospital. "This is an extremely rare occurrence," said an unnamed source at the scene. A second group of tourists were on the bus at the time of the incident but were not injured, according to reports.

GPT-Neo



---

Extracted Features:

$\log p(x)$  (GPT-2): 3.43  
 $\log p(x)$  (GPT-Neo):3.21  
 $\log p(x)$  (GPT-J): 3.23  
 $\log p(x)$  (LLaMA): 3.08

pct-scores / Pearson scores / Spearman scores:

GPT-2 and GPT-Neo	0.54/0.96/0.95
GPT-2 and GPT-J	0.40/0.95/0.95
GPT-2 and LLaMA	0.56/0.83/0.87
GPT-Neo and GPT-J	0.42/0.97/0.98
GPT-Neo and LLaMA	0.49/0.86/0.91
GPT-J and LLaMA	0.56/0.89/0.93

---

Tracing Results:

GPT-2: 0.3%  
GPT-J/Neo: 72.6 %  
LLaMA: 3.0%  
ChatGPT: 18.3 %  
Human: 5.8 %

---

Table 7: Random selected case studies (a).

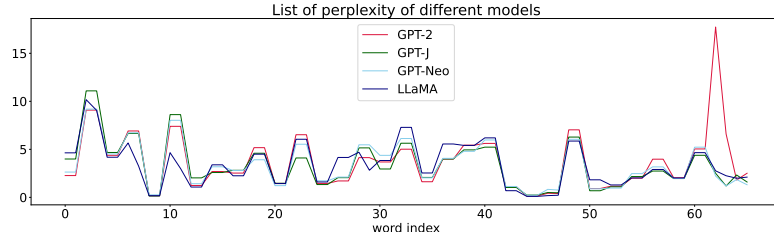
---

Text Origin

---

He said ministers who disagreed with austerity measures could not go to other countries for talks as Germany would only make life difficult for them. And he said that he was a “proud German”.

GPT-J



---

Extracted Features:

$\log p(x)$  (GPT-2): 3.86  
 $\log p(x)$  (GPT-Neo): 2.99  
 $\log p(x)$  (GPT-J): 2.96  
 $\log p(x)$  (LLaMA): 3.07

pct-scores / Pearson scores / Spearman scores:

GPT-2 and GPT-Neo	0.65/0.70/0.90
GPT-2 and GPT-J	0.47/0.70/0.89
GPT-2 and LLaMA	0.44/0.65/0.85
GPT-Neo and GPT-J	0.44/0.96/0.98
GPT-Neo and LLaMA	0.50/0.86/0.90
GPT-J and LLaMA	0.56/0.85/0.89

---

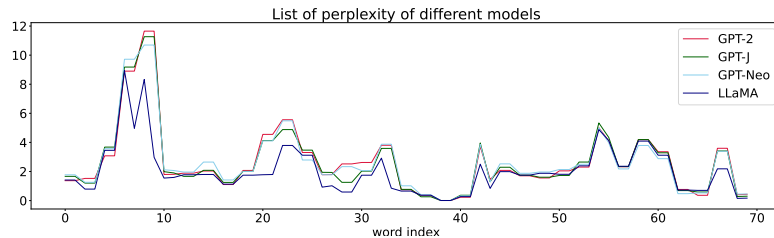
Tracing Results:

GPT-2: 0.0%  
GPT-J/Neo: 100.0 %  
LLaMA: 0.0%  
ChatGPT: 0.0 %  
Human: 0.0 %

---

This is a very memorable spaghetti western. It has a great storyline about two brothers, one a sheriff and the other a gunfighter. It is a very good one. I will recommend it to my friends.

LLaMA



---

Extracted Features:

$\log p(x)$  (GPT-2): 2.48  
 $\log p(x)$  (GPT-Neo): 2.46  
 $\log p(x)$  (GPT-J): 2.42  
 $\log p(x)$  (LLaMA): 1.91

pct-scores / Pearson scores / Spearman scores:

GPT-2 and GPT-Neo	0.57/0.99/0.97
GPT-2 and GPT-J	0.56/0.99/0.95
GPT-2 and LLaMA	0.34/0.83/0.85
GPT-Neo and GPT-J	0.47/0.99/0.96
GPT-Neo and LLaMA	0.21/0.85/0.87
GPT-J and LLaMA	0.23/0.86/0.91

---

Tracing Results:

GPT-2: 0.2 %  
GPT-J/Neo: 0.3 %  
LLaMA: 59.1 %  
ChatGPT: 31.3 %  
Human: 9.0 %

---

Table 8: Random selected case studies (b).

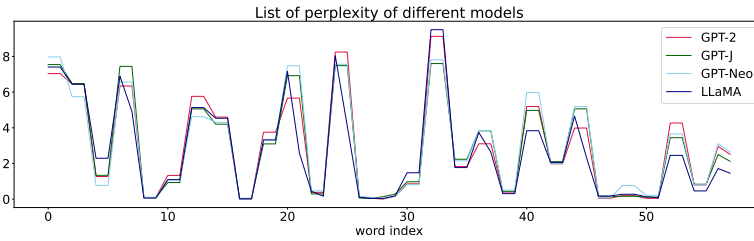
---

Text Origin

---

The film in question focuses on the character played by the immensely talented Helena Bonham Carter, who delivered an impressive performance despite being confined to a wheelchair throughout the movie.

Chat-GPT



---

Extracted Features:

$\log p(x)$  (GPT-2): 2.73  
 $\log p(x)$  (GPT-Neo): 2.77  
 $\log p(x)$  (GPT-J): 2.72  
 $\log p(x)$  (LLaMA): 2.19

pct-scores / Pearson scores / Spearman scores:

GPT-2 and GPT-Neo	0.64/0.97/0.98
GPT-2 and GPT-J	0.55/0.98/0.98
GPT-2 and LLaMA	0.47/0.94/0.97
GPT-Neo and GPT-J	0.38/0.99/0.98
GPT-Neo and LLaMA	0.41/0.91/0.95
GPT-J and LLaMA	0.41/0.93/0.97

---

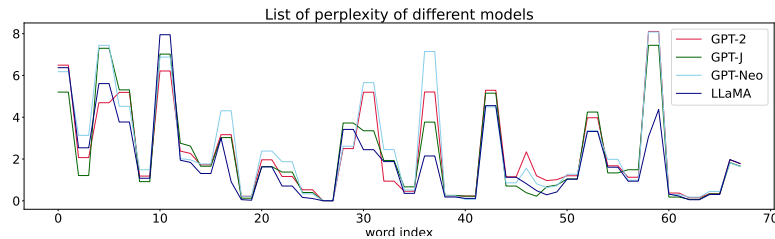
Tracing Results:

GPT-2: 0.8%  
GPT-J/Neo: 0.3 %  
LLaMA: 0.6 %  
ChatGPT: 97.3 %  
Human: 1.0 %

---

College sports are also popular in southern California. The UCLA Bruins and the USC Trojans both field teams in NCAA Division I in the Pac-12 Conference, and there is a longtime rivalry between the schools.

Human



---

Extracted Features:

$\log p(x)$  (GPT-2): 2.03  
 $\log p(x)$  (GPT-Neo): 2.21  
 $\log p(x)$  (GPT-J): 1.95  
 $\log p(x)$  (LLaMA): 1.74

pct-scores / Pearson scores / Spearman scores:

GPT-2 and GPT-Neo	0.59/0.95/0.95
GPT-2 and GPT-J	0.35/0.93/0.92
GPT-2 and LLaMA	0.21/0.85/0.92
GPT-Neo and GPT-J	0.41/0.93/0.95
GPT-Neo and LLaMA	0.22/0.84/0.94
GPT-J and LLaMA	0.26/0.91/0.94

---

Tracing Results:

GPT-2: 0.0 %  
GPT-J/Neo: 0.0 %  
LLaMA: 14.1%  
ChatGPT: 40.6 %  
Human: 45.3 %

---

Table 9: Random selected case studies (c).



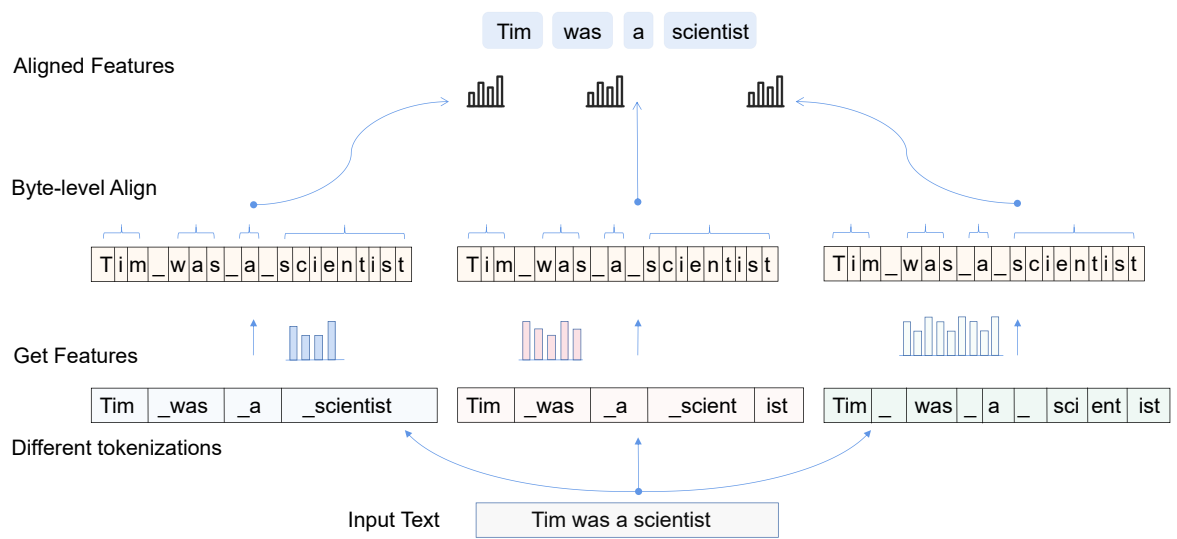


Figure 6: Align Strategies